

(19) World Intellectual Property  
Organization  
International Bureau



(43) International Publication Date  
27 January 2005 (27.01.2005)

PCT

(10) International Publication Number  
**WO 2005/008517 A1**

(51) International Patent Classification<sup>7</sup>: **G06F 17/18**

(21) International Application Number:  
PCT/AU2003/000923

(22) International Filing Date: 18 July 2003 (18.07.2003)

(25) Filing Language: English

(26) Publication Language: English

(71) Applicant (for all designated States except US): **COMMONWEALTH SCIENTIFIC AND INDUSTRIAL RESEARCH ORGANISATION** [AU/AU]; Limestone Avenue, CAMPBELL, Australian Capital Territory 2612 (AU).

(72) Inventor; and

(75) Inventor/Applicant (for US only): **STONE, Glenn** [AU/AU]; 19/55-61 Old Northern Road, BAULKHAM HILLS, New South Wales 2153 (AU).

(74) Agent: **GRIFFITH HACK**; GPO Box 4164, SYDNEY, New South Wales 2001 (AU).

(81) Designated States (*national*): AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DE, DK, DM, DZ, EC, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ, NI, NO, NZ, OM, PG, PH, PL, PT, RO, RU, SC, SD, SE, SG, SK, SL, SY, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, YU, ZA, ZM, ZW.

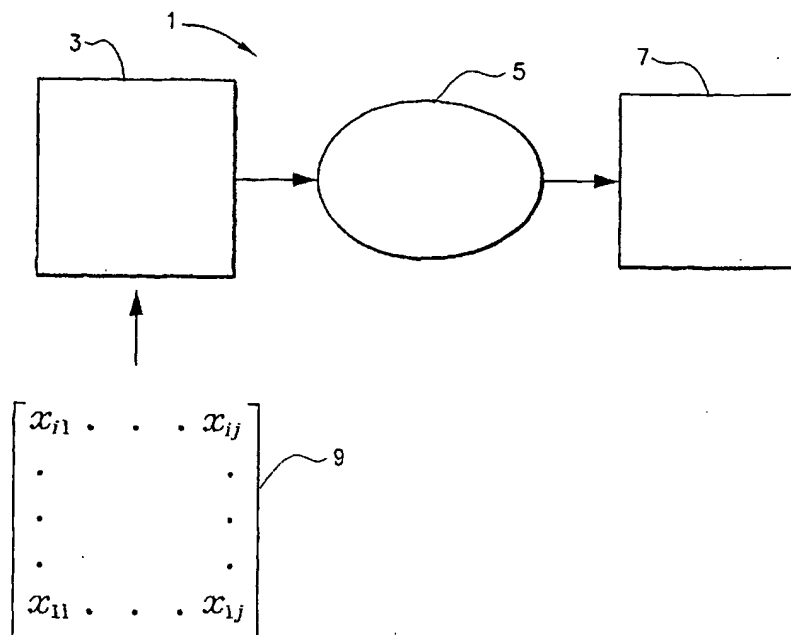
(84) Designated States (*regional*): ARIPO patent (GH, GM, KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZM, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HU, IE, IT, LU, MC, NL, PT, RO, SE, SI, SK, TR), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG).

Published:

— with international search report

For two-letter codes and other abbreviations, refer to the "Guidance Notes on Codes and Abbreviations" appearing at the beginning of each regular issue of the PCT Gazette.

(54) Title: A METHOD AND SYSTEM FOR SELECTING ONE OR MORE VARIABLES FOR USE WITH A STATISTICAL MODEL



(57) Abstract: A method of selecting one or more variables for use with a statistical model, the method comprising the steps of: creating a plurality of unique subsets of variables of multivariate data; determining the performance of a discriminant rule when used with each of the subsets, the discriminant rule being based on multivariate normal class densities each having substantially diagonal covariance matrices; and selecting the one or more variables from at least one of the subsets that result in a desired performance of the discriminant rule.



WO 2005/008517 A1

- 1 -

A METHOD AND SYSTEM FOR SELECTING ONE OR MORE VARIABLES FOR  
USE WITH A STATISTICAL MODEL

FIELD OF THE INVENTION

5

The present invention relates to a system and method for selecting one or more variables for use with a statistical model. The present invention is of particular, but by no means exclusive, application to building a  
10 classifier that is capable of predicting the class of an observation.

BACKGROUND OF THE INVENTION

15

Generally speaking, a statistical model is a description of an assumed structure of a set of observations. Typically, the statistical model is in the form of a mathematical function of the process assumed to have generated the observations. The mathematical function  
20 is usually dependent on a number of variables that have been carefully selected to ensure the mathematical function accurately models the assumed process.

SUMMARY OF THE INVENTION

25

According to a first aspect of the present invention, there is provided a method of selecting one or more variables for use with a statistical model, the method comprising the steps of:

30

creating a plurality of unique subsets of variables of multivariate data;

determining the performance of a discriminant rule when used with each of the subsets, the discriminant rule being based on multivariate normal class densities  
35 each having substantially diagonal covariance matrices; and selecting the one or more variables from at least one of the subsets that result in a desired performance of

- 2 -

the discriminant rule.

Given that the discriminant rule used in the method is widely considered to be suitable only for independent multinormal data, studies by the applicant have surprising shown that that method is in fact well suited to some data that is not independent multinormal, for example gene expression data.

Preferably, the step of creating the plurality of unique subsets comprises the step of identifying a variable in the multivariate data that is not a member of a set of variables, and adding the identified variable to the set.

This approach to creating the subsets is based on a forward stepwise variable selection technique.

Alternatively, the step of creating the plurality of unique subsets comprises the step of identifying a variable in the set which has not been previously removed, and removing the identified variable from the set.

This alternative approach is based on a backward stepwise variable selection technique.

Preferably, the step of determining the performance of the discriminant rule comprises assessing a prediction error rate of the discriminant rule.

Even more preferably, the prediction error rate is a cross-validated error rate.

Alternatively, the step of determining the performance of the discriminant rule is assessed using a likelihood based approach.

Preferably, the desired performance of the

discriminant rule comprises the lowest possible prediction error rate of the discriminant rule.

Alternatively, the desired performance may be any other desired error rate.

Preferably, the multivariate data comprises gene expression data.

According to a second aspect of the present invention, there is provided computer software which, when executed by a computer, enables the computer to carry out the steps described in the first aspect of the present invention.

According to a third aspect of the present invention, there is provided a computer storage medium containing the software described in the second aspect of the present invention.

According to a fourth aspect of the present invention, there is provided a statistical model for predicting a class of an observation, wherein the model includes one or more variables that have been selected using the method described in the first aspect of the present invention.

According to a fifth aspect of the present invention, there is provided an apparatus for selecting one or more variables for use with a statistical model, the system comprising:

data creating means arranged to create a plurality of unique subsets of variables of multivariate data;

a processing means arranged to determine the performance of a discriminant rule when used with each of the subsets, the discriminant rule being based on

- 4 -

multivariate normal class densities each having substantially diagonal covariance matrices; and

5 a selecting means arranged to select the one or more variables from at least one of the subsets that results in a desired performance of the discriminant rule.

10 Preferably, the data creating means is arranged to create the plurality of unique subsets by identifying a variable in the multivariate data that is not a member of a set of variables, and adding the identified variable to the set.

15 Alternatively, the data creating means is arranged to create the plurality of unique subsets by identifying a variable in the set which has not been previously removed, and removing the identified variable from the set.

20 Preferably, the determining means is arranged to determine the performance of the discriminant rule by assessing a prediction error rate of the discriminant rule.

25 Even more preferably, the prediction error rate is a cross-validated error rate.

30 Alternatively, the determining means is arranged to determine the performance of the discriminant rule using a likelihood based approach.

35 Preferably, the desired performance of the discriminant rule comprises the lowest possible prediction error rate of the discriminant rule.

Alternatively, the desired performance may be any other desired error rate.

- 5 -

Preferably, the multivariate data comprises gene expression data.

Preferably, the data creating means, processing means and selecting means are in the form of a computer running software.

#### BRIEF DESCRIPTION OF THE DRAWINGS

Notwithstanding any other embodiments that may fall within the scope of the present invention, a preferred embodiment of the present invention will now be described, by way of example only, with reference to the accompanying figures, in which:

Figure 1, illustrates a block diagram of the components that are included in an apparatus, according to the preferred embodiment of the present invention, that is arranged to select one or more variables for use with a statistical model; and

Figure 2 illustrates a flow diagram of the various steps carried out by the apparatus of figure 1.

#### A PREFERRED EMBODIMENT OF THE INVENTION

As can be seen in figure 1, an apparatus 1 according to the preferred embodiment of the present invention comprises data creating means 3, processing means 5, and selecting means 7. The data creating means 3, processing means 5 and selecting means 7 are in the form of a computer running software.

The data creating means 3 is arranged such that it has access to multivariate data 9; that is data for which each observation consists of values for more than one variable. In the preferred embodiment the multivariate data is gene expression data. An example of gene expression data is the leukemia data set referred to in the article

entitled "Molecular classification of cancer: class discovery and class prediction by gene expression monitoring", which appeared in *Science* 286:531-537, 1999.

5           The data creating means 3 processes the multivariate data 9 in order to produce a plurality of unique subsets of variables of the multivariate data 9.

Essentially, the data creating means 3 creates the  
10 plurality of unique subsets by employing a technique that is similar to forward stepwise variable selection. Generally speaking, forward stepwise selection involves identifying those variables in the multivariate data that are not in a set of variables which are 'in a statistical  
15 model', and adding them to the set one at a time. It is the process of adding the variables to the set that results in the creations of the plurality of unique subsets. Further details on the forward stepwise variable selection technique can be found in most texts covering discriminant  
20 function analysis. One such text can be found on the Internet at  
<http://www.statsoftinc.com/textbook/stdiscan.html>

Following the addition of a variable to the set,  
25 the processing means 5 applies the set (which is effectively one of the plurality of unique subsets) to a discriminant rule, and makes a record of the performance of the discriminant rule when used with the variables in the set. The processing means 5 continues this processes for  
30 each variable added to the set; that is, the processing means records the performance of the discriminant rule for each one of the unique subsets.

The discriminant rule used by the processing  
35 means 5 is based on multivariate normal class densities each having substantially diagonal covariance matrices, and is in the form of one of the following functions:

- 7 -

$$C(x) = \operatorname{argmin}_k \sum_{j=1}^n \left\{ \frac{(x_j - \mu_{kj})^2}{\sigma_{kj}^2} + \log \sigma_{kj}^2 \right\} \quad (1)$$

$$C(x) = \operatorname{argmin}_k \sum_{j=1}^n \frac{(x_j - \mu_{kj})^2}{\sigma_j^2} \quad (2)$$

5           The first function (1) assumes that the class densities have diagonal covariance matrices,  $\Delta_k = \operatorname{diag}(\sigma_{k1}^2, \dots, \sigma_{kp}^2)$ , whilst the second function (2) assumes the class densities have the same diagonal covariance matrix,  $\Delta_k = \operatorname{diag}(\sigma_1^2, \dots, \sigma_p^2)$ .

10

In order to determine the performance of the discriminant rule, the processing means 5 is arranged to determine the cross-validated error rate of the predictor.

15

Once the processing means 5 has applied each of the unique subsets to the discriminant rule, the processing means 5 examines the recorded error rates to identify the subset that results in the lowest error rate. The processing means 5 then proceeds to select the one or more variables (for use with the statistical model) from the identified subset (that is, the subset that results in the lowest error rate) as the variables to be used with the statistical model.

20

25

The use of the forward stepwise technique means that the apparatus 1 is effectively performing the following steps:

30

1. Starting with an empty set of variables;
2. For each variable of the multivariate data not in the set, add to set and determine the performance of the discriminant rule;
3. Add variable to the set which results in the discriminant rule having the best performance; and

35



4. Continuing steps 1 - 3 while the performance of the discriminant rule is improving.

In order to select the one or more variables for use with the statistical model, the apparatus 1 is effectively carrying out the following broad steps:

creating a plurality of unique subsets of variables of multivariate data;

determining the performance of the discriminant rule when used with each of the subsets, the discriminant rule being based on multivariate normal class densities each having substantially diagonal covariance matrices; and

selecting the one or more variables from at least one of the subsets that result in a desired performance of the discriminant rule.

In order to gain an insight into the performance of the preferred embodiment of the present invention, the preferred embodiment was applied to Alizadeh's DLBCL data. The DLBCL data can be obtained from <http://genome-www.stanfordd.edu/lymphoma>. This data was collected from 42 patients and represents two classes of diffuse large B-cell lymphoma (DLBCL), GC and Activated. The preferred embodiment of the present invention selected just three genes (variables) from the DLBCL data. The three genes were then used in a classification which produced no errors (re-substitution), and when cross-validated the classifier produced about 5 errors (approximately 12%).

It is noted that whilst the preferred embodiment uses the cross-validated error rate as a measure of the discriminant rule's performance, other techniques for determining the performance of the discriminant rule are considered to be suitable. For example, a likelihood based approach.

Whilst the preferred embodiment employs a forward

- 9 -

stepwise variable selection technique to create the plurality of unique subsets, it is envisaged that alternative techniques such a backward stepwise variable selection could be used with the present invention.

5

It will be appreciated that whilst the description of the preferred embodiment refers to the multivariate data as being gene expression data, the present invention can be used with multivariate data other  
10 that gene expression data.

Those skilled in the art will appreciate that the invention described herein is susceptible to variations and modifications other than those specifically described. It  
15 should be understood that the invention includes all such variations and modifications which fall within the spirit and scope of the invention.

- 10 -

## CLAIMS:

1. A method of selecting one or more variables for use with a statistical model, the method comprising the  
5 steps of:

creating a plurality of unique subsets of variables of multivariate data;

determining the performance of a discriminant rule when used with each of the subsets, the discriminant  
10 rule being based on multivariate normal class densities each having substantially diagonal covariance matrices; and

selecting the one or more variables from at least one of the subsets that result in a desired performance of the discriminant rule.

15

2. The method as claimed in claim 1, wherein the step of creating the plurality of unique subsets comprises the step of identifying a variable in the multivariate data that is not a member of a set of  
20 variables, and adding the identified variable to the set.

3. The method as claimed in any one of claims 1 or 2, wherein the step of determining the performance of the discriminant rule comprises assessing a prediction  
25 error rate of the discriminant rule.

4. The method as claimed in claim 3, wherein the prediction error rate is a cross-validated error rate.

5. The method as claimed in any one of the preceding claims, wherein the desired performance of the discriminant rule comprises the lowest possible prediction error rate of the discriminant rule.

6. The method as claimed in any one of the preceding claims, wherein the multivariate data comprises gene expression data.

- 11 -

7. Computer software which, when executed by a computer, enables the computer to carry out the steps defined in any one of the preceding steps.

5 8. A computer storage medium containing the software defined in claim 7.

9. A statistical model for predicting a class of an observation, wherein the model includes one or more  
10 variables that have been selected using the method defined in any one of claims 1 - 6.

10. An apparatus for selecting one or more variables for use with a statistical model, the system  
15 comprising:

data creating means arranged to create a plurality of unique subsets of variables of multivariate data;

a processing means arranged to determine the  
20 performance of a discriminant rule when used with each of the subsets, the discriminant rule being based on multivariate normal class densities each having substantially diagonal covariance matrices; and

a selecting means arranged to select the one or  
25 more variables from at least one of the subsets that results in a desired performance of the discriminant rule.

11. The apparatus as claimed in claim 10, wherein the data creating means is arranged to create the  
30 plurality of unique subsets by identifying a variable in the multivariate data that is not a member of a set of variables, and adding the identified variable to the set.

12. The apparatus as claimed in any one of  
35 claims 10 or 11, wherein the determining means is arranged to determine the performance of the discriminant rule by assessing a prediction error rate of the discriminant rule.

- 12 -

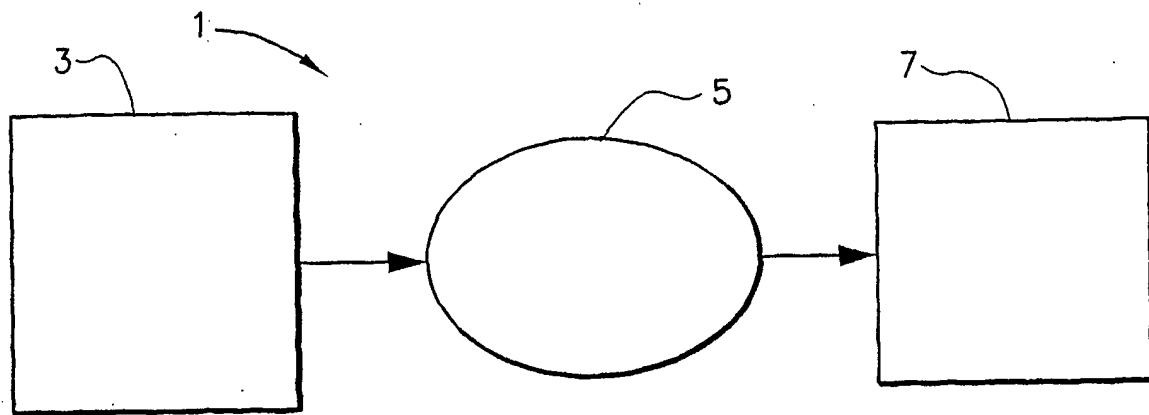
13. The apparatus as claimed in claim 12, wherein the prediction error rate is a cross-validated error rate.

14. The apparatus as claimed in any one of the  
5 preceding claims, wherein the desired performance of the discriminant rule comprises the lowest possible prediction error rate of the discriminant rule.

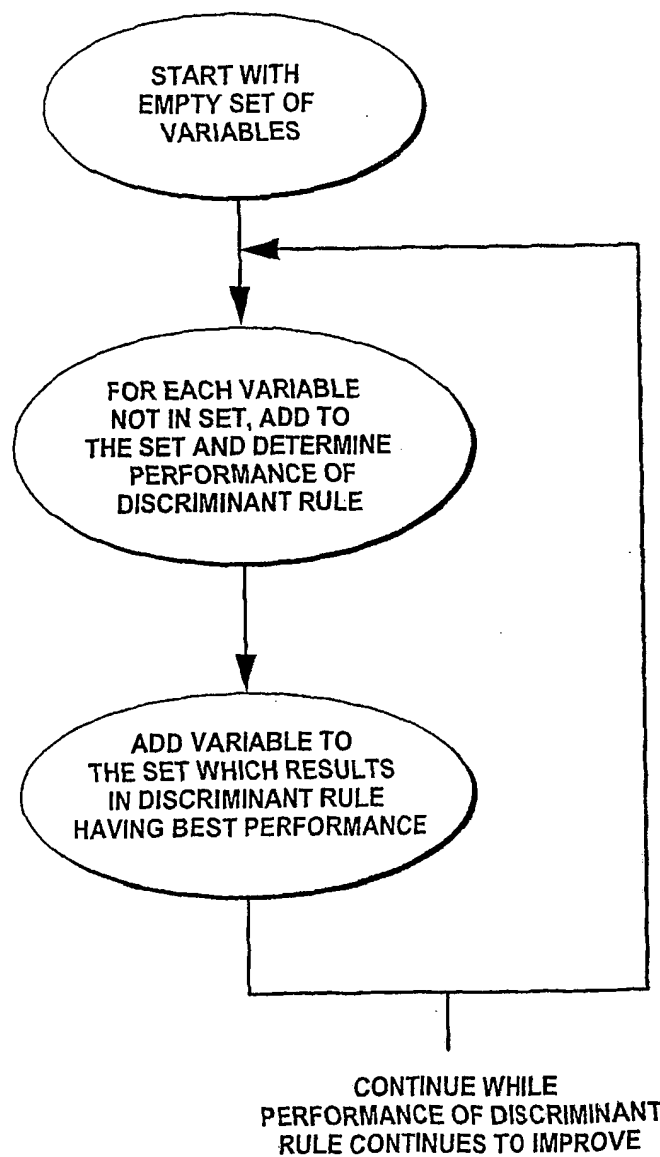
15. The apparatus as claimed in any one of  
10 claims 10 - 14, wherein the multivariate data comprises gene expression data.

16. The apparatus as claimed in any one of  
claims 10 - 15, wherein the data creating means, processing  
15 means and selecting means are in the form of a computer running software.

1/1

**Fig. 1**

$$\begin{bmatrix} x_{i1} & \cdot & \cdot & \cdot & x_{ij} \\ \cdot & & & & \cdot \\ \cdot & & & & \cdot \\ \cdot & & & & \cdot \\ x_{11} & \cdot & \cdot & \cdot & x_{1j} \end{bmatrix} \quad 9$$

**Fig. 2**

# INTERNATIONAL SEARCH REPORT

International application No.

PCT/AU03/00923

## A. CLASSIFICATION OF SUBJECT MATTER

Int. Cl. <sup>7</sup>: G06F 17/18

According to International Patent Classification (IPC) or to both national classification and IPC

## B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)  
WPAT: IPC and keywords. Keywords included classif+, discriminant and statistic.

## C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	US 5970239 A (Bahl et al) 19 October 1999 - whole document	1 to 16
X	WO 98/32088 A (Chiron Corporation) 23 July 1998 - whole document	1 to 16
X	WO 02/25405 A2 (The Regents of the University of California) 28 March 2002 - whole document	1 to 16

☒ Further documents are listed in the continuation of Box C

☒ See patent family annex

- \* Special categories of cited documents:
- |   |  |
|---|--|
| "A" document defining the general state of the art which is not considered to be of particular relevance  | "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention  |
| "E" earlier application or patent but published on or after the international filing date   | "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone   |
| "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) | "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art |
| "O" document referring to an oral disclosure, use, exhibition or other means  | "&" document member of the same patent family  |
| "P" document published prior to the international filing date but later than the priority date claimed  |  |

Date of the actual completion of the international search 8 August 2003	Date of mailing of the international search report 14 AUG 2003
Name and mailing address of the ISA/AU AUSTRALIAN PATENT OFFICE PO BOX 200, WODEN ACT 2606, AUSTRALIA E-mail address: pct@ipaustalia.gov.au Facsimile No. (02) 6285 3929	Authorized officer  J.W. THOMSON Telephone No : (02) 6283 2214

## INTERNATIONAL SEARCH REPORT

International application No.

PCT/AU03/00923

C (Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	Dr Min Qiu, 'Multivariate Discriminant Analysis' Advanced Data Analysis, Information Management and Marketing, University of Western Australia 4 August 2002 [retrieved on 11 August 2003] Retrieved from the Internet: URL: <a href="http://www.imm.ece.uwa.edu.au/unit450461/lectures/450461_week5.pdf">http://www.imm.ece.uwa.edu.au/unit450461/lectures/450461_week5.pdf</a>	1 to 16
A	EP 501784 B1 (Philip Morris Products Inc) 2 September 1992 (note column 5 line 44 to column 12 line 33)	



# INTERNATIONAL SEARCH REPORT

Information on patent family members

International application No.

PCT/AU03/00923

This Annex lists the known "A" publication level patent family members relating to the patent documents cited in the above-mentioned international search report. The Australian Patent Office is in no way liable for these particulars which are merely given for the purpose of information.

Patent Document Cited in Search Report				Patent Family Member			
US	5970239	NONE					
WO	200225405	AU	200194644	US	2002111742		
EP	501784	US	5146510	CA	2061865	HK	1013872
		JP	5126757	US	5353356	EP	382466
		JP	2242482	US	5046111	US	5165101
		US	5189708				
WO	9832088	AU	60237/98	EP	953177	US	5860917
END OF ANNEX							